

$$\vec{\epsilon} \sim N_n(\vec{0}_n, \sigma^2 I_p), \vec{Y} = X\vec{\beta} + \vec{\epsilon}, \vec{\beta} \sim N_{p+1}(\vec{0}_{p+1}, \tau^2 I_{p+1}), f(\sigma^2) \propto \frac{1}{\sigma^2}$$

$$\vec{\beta}_{\text{Ridge}} \Rightarrow \vec{\beta}^{\text{MAP}} = (X^T X + \lambda I_{p+1})^{-1} X^T \vec{Y} \quad \text{where } \lambda = \frac{\sigma^2}{\tau^2} > 0$$

Ridge Estimate

If we employ this estimator for prediction purposes, we ignore the fact that σ^2 is unknown and hence λ as a hyperparameter

Slighly processed problem: we shrink each slope estimate equally as σ^2 is the same prior variance. If the \vec{x}_{ij} 's are on different scales, it is doing unequal shrinking. So, as a prestep, normalize all variables:

let X be the design matrix where $\vec{x}_{ij} = \frac{\bar{x}_{ij} - \bar{x}_{\cdot j}}{s_{\vec{x}_{ij}}}$ for $j \geq 1$
(omit intercept)

this is called "standardization". The standardized feature responses have avg zero, std. err one. They all shrink equally. Demo

So assuming standardized covariances...

6

let's say prior on $\vec{\beta}$:

$$\beta_0, \beta_1, \dots, \beta_p \stackrel{\text{iid}}{\sim} \text{Laplace}(0, \sigma^2) := \frac{1}{2\sigma^2} e^{-\frac{|\beta_j|}{\sigma^2}}$$

$$f(\sigma^2) \propto \frac{1}{\sigma^2}$$

$$f(\vec{\beta}, \sigma^2 | X, \vec{y}) \propto f(\vec{y} | \vec{\beta}, \sigma^2, X) f(\vec{\beta}, \sigma^2 | X) \\ = \frac{1}{(2\pi)^{n/2} \sqrt{\det(\sigma^2 I_n)}} e^{-\frac{1}{2}(\vec{y} - X\vec{\beta})^\top (\sigma^2 I_n)^{-1} (\vec{y} - X\vec{\beta})}$$

$$\prod_{j=0}^p \frac{1}{2\sigma^2} e^{-\frac{|\beta_j|}{\sigma^2}} \cdot \frac{1}{\sigma^2}$$

$$\propto (\sigma^2)^{-\frac{n}{2}-1} e^{-\frac{1}{2\sigma^2} \|\vec{y} - X\vec{\beta}\|^2} e^{-\frac{1}{\sigma^2} \sum_{j=0}^p |\beta_j|}$$

Now find $\vec{\beta}^{\text{map}}$

$$\vec{\beta}^{\text{map}} = \underset{\vec{\beta}}{\text{argmax}} \{ f(\vec{\beta}, \sigma^2 | X, \vec{y}) \} = \underset{\vec{\beta}}{\text{argmax}} \{ \ln f(\vec{\beta}, \sigma^2 | X, \vec{y}) \}$$

$$= \underset{\vec{\beta}}{\text{argmax}} \left\{ \underbrace{\left(-\frac{n}{2}-1\right) \ln \sigma^2}_{\text{constant}} - \frac{1}{2\sigma^2} \|\vec{y} - X\vec{\beta}\|^2 - \frac{1}{\sigma^2} \sum_{j=0}^p |\beta_j| \right\}$$

$$= \underset{\vec{\beta}}{\text{argmax}} \left\{ \underbrace{-\frac{1}{2\sigma^2} \left(\|\vec{y} - X\vec{\beta}\|^2 + \frac{2\sigma^2}{\sigma^2} \sum_{j=0}^p |\beta_j| \right)}_{\text{constant}} \right\}$$

$$= \underset{\vec{\beta}}{\text{argmin}} \left\{ \underbrace{\|\vec{y} - X\vec{\beta}\|^2}_{\text{minimize}} + \frac{2\sigma^2}{\sigma^2} \sum_{j=0}^p |\beta_j| \right\}$$

Now join the 3rd term...

$$\text{let } \lambda = \frac{2\sigma^2}{\sigma^2}$$

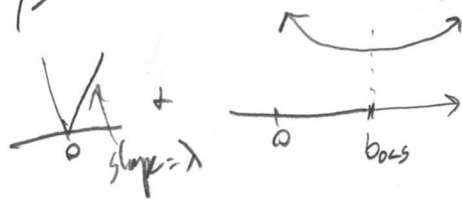
$$\vec{b}_{\text{LASSO}} := \underset{\vec{\beta}}{\text{argmin}} \left\{ \text{SSE} + \lambda \sum_{j=0}^p |\beta_j| \right\}$$

No closed form solution exists.

Need to use optimizer

As $\lambda \rightarrow 0 \Rightarrow \vec{b}_{\text{Lasso}} \rightarrow \vec{b}$. As $\lambda \rightarrow \infty \Rightarrow \vec{b}_{\text{Lasso}} \rightarrow \vec{0}_{p \times 1}$

Also. If $\vec{\beta}$ was one dimensional, here's what the lasso problem looks like (demo)



If λ is large and/or if $\beta_{\text{lasso}} \neq 0$, the sum will have min. near zero (demo)

It's difficult for the $\|y - X\beta\|^2$ term to "beat" the $\lambda \sum_{j=1}^p |\beta_j|$ term. Hence, Lasso has an argument at $\beta = 0$ due to the sharp L1 norm which has its sharp mode at zero.

\Rightarrow Lasso has an ability to do "variable selection".

I.e. it ^{thins} picks variables when there is a large set that should be pared down as you know you'll overfit if you use all of them